

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problems Mailbox.**

(19)



JAPANESE PATENT OFFICE

PATENT ABSTRACTS OF JAPAN

(11) Publication number: **10260991 A**

(43) Date of publication of application: **29.09.98**

(51) Int. Cl **G06F 17/30**

(21) Application number: **09270251**

(22) Date of filing: **02.10.97**

(30) Priority: **14.01.97 JP 09 5010**

(71) Applicant: **SEIKO EPSON CORP**

(72) Inventor: **MIWA SHINJI**

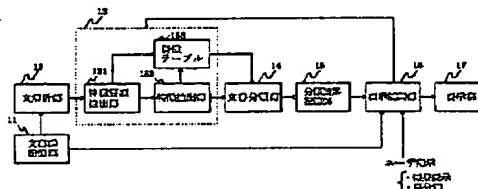
**(54) INFORMATION RETRIEVING METHOD AND
INFORMATION RETRIEVING DEVICE**

(57) Abstract:

PROBLEM TO BE SOLVED: To accurately take out information that is needed for a user by watching the characteristic element of each displayed cluster to select a cluster and performing an operation, that interactively and also gradually narrows a retrieval object, of electing either resorting or result display.

SOLUTION: A characteristic extracting part 132 counts how many times which characteristic element appears based on characteristic elements that are extracted from a characteristic element extracting part 131 and generates a characteristic table 133. A document sorting part 14 refers to the table 133 and divides into plural clusters based on statistical information such as the appearance frequency of each characteristic element that exists in each document. A display controlling part 16 controls that the content of a sort result storing part 15 is shown as a sort result on a display part 17, also forms display data as a retrieval result content based on the content of the part 15 and the content of a document group storing part 11 when an instruction of retrieval result display is made by a user and shows the display data on the part 17.

COPYRIGHT: (C)1998,JPO



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平10-260991

(43) 公開日 平成10年(1998) 9月29日

(51) Int.Cl.⁹

G 0 6 F 17/30

識別記号

F I

G 0 6 F 15/401

15/40

3 1 0 D

3 7 0 A

審査請求 未請求 請求項の数6 O L (全 9 頁)

(21) 出願番号 特願平9-270251

(22) 出願日 平成9年(1997)10月2日

(31) 優先権主張番号 特願平9-5010

(32) 優先日 平9(1997)1月14日

(33) 優先権主張国 日本 (J P)

(71) 出願人 000002369

セイコーエプソン株式会社

東京都新宿区西新宿2丁目4番1号

(72) 発明者 三輪 真司

長野県諏訪市大和3丁目3番5号 セイコーエプソン株式会社内

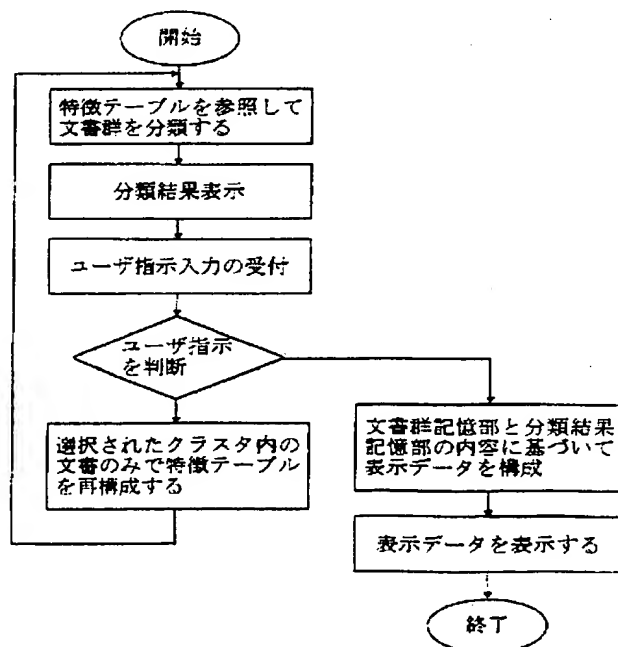
(74) 代理人 弁理士 鈴木 喜三郎 (外2名)

(54) 【発明の名称】 情報検索方法および情報検索装置

(57) 【要約】

【課題】 適当なキーワードの指定が難しい場合でも的確な情報検索を可能とする。

【解決手段】 それぞれの文書から特徴要素を抽出し、その特徴要素とその特徴要素を含む文書との関係を表す特徴テーブルを作成して、その特徴テーブルを用いて文書群を複数のクラスタに分類して表示する(ステップs1, s2)。このクラスタは、各クラスタごとにそのクラスタを代表する特徴要素とその特徴要素を一定以上含む文書数とをデータとして有する。そして、ユーザからのクラスタ選択指示を受けたとき、選択されたクラスタに属する文書内容の表示指示または再分類指示を受け付けて、内容表示指示の場合は、当該クラスタに属する文書内容の表示を行う(ステップs3~s6)。一方、再分類指示の場合は、当該クラスタに属する文書のみで特徴テーブルを再構成して(ステップs3~s5)、その再構成された特徴テーブルに基づいてクラスタに分類して表示する。



【特許請求の範囲】

【請求項1】 文書群に属するそれぞれの文書を解析し、それぞれの文書から特徴要素を抽出し、その特徴要素とその特徴要素を含む文書との関係を表す特徴テーブルを作成して、その特徴テーブルに基づいて文書群を複数のクラスタに分類して表示し、ユーザからのクラスタ選択指示を受けたとき、その選択されたクラスタに属する文書に関する内容の表示指示または再分類指示を受け付けて、内容表示指示の場合は、当該クラスタに属する文書に関する表示を行い、再分類指示の場合は、当該クラスタに属する文書のみで前記特徴テーブルを再構成してその再構成された特徴テーブルに基づいてクラスタに分類して表示することを特徴とする情報検索方法。

【請求項2】 前記特徴テーブルに基づいて文書群を複数のクラスタに分類する処理は、それぞれの文書内に存在するそれぞれの特徴要素の出現頻度などの統計的な情報に基づいて複数のクラスタに分類することを特徴とする請求項1記載の情報検索方法。

【請求項3】 前記分類された複数のクラスタは、少なくとも、それぞれのクラスタごとにそのクラスタを代表する特徴要素と、その特徴要素を一定以上含む文書数とをデータとして有することを特徴とする請求項2記載の情報検索方法。

【請求項4】 文書群を記憶する文書群記憶部と、この文書群記憶部に記憶されているそれぞれの文書を解析する文解析部と、この文解析部による解析結果からそれぞれの文書に対する特徴要素を抽出し、その特徴要素とその特徴要素を含む文書との関係を表す特徴テーブルを作成する特徴テーブル作成部と、前記特徴テーブルの内容に基づいて文書群を複数のクラスタに分類する文書分類部と、この文書分類部により分類された内容を記憶する分類結果記憶部と、この分類結果記憶部の内容を読み出して複数のクラスタを表示させる制御を行うとともに、ユーザからのクラスタ選択指示を受けたとき、その選択されたクラスタに属する文書に関する内容の表示指示または再分類指示を受け付けて、内容表示指示の場合は、当該クラスタに属する文書を表示させる制御を行い、再分類指示の場合は、当該クラスタに属する文書のみで前記特徴テーブルを再構成させる制御を行う表示制御部と、を有することを特徴とする情報検索装置。

【請求項5】 前記情報分類部が行う特徴テーブルに基づいてそれぞれの文書を複数のクラスタに分類する処理は、それぞれの文書内に存在するそれぞれの特徴要素の出現頻度などの統計的な情報に基づいて複数のクラスタに分類することを特徴とする請求項4記載の情報検索装置。

【請求項6】 前記分類された複数のクラスタは、少な

くとも、それぞれのクラスタごとにそのクラスタを代表する特徴要素と、その特徴要素を一定以上含む文書数とをデータとして有することを特徴とする請求項5記載の情報検索装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、文書データを蓄積したデータベースやネットワークで公開されている文書群から、ユーザの要求に合致する文書を検索し、提示し得る情報検索方法および情報検索装置に関する。

【0002】

【従来の技術】文書データを蓄積したデータベースなどからユーザの要求する文書を効率よく検索して取り出す方法としては、ユーザの入力したキーワードをもとに文書データを取り出すというような検索方法が一般的である。

【0003】この検索方法は、検索しようとする文書に対してユーザ自身が何らかのキーワードを考えて、そのキーワードを入力することにより、システム側で、そのキーワードに合致する文書を取り出して出力するというものである。

【0004】

【発明が解決しようとする課題】しかしながら、上述したような情報検索方法は、検索対象のデータベースの規模や格納されている文書の種類にかかわらず、入力されたキーワードのみによる検索であるため、状況によっては、検索されて出力される文書量が膨大な量となったり、逆に、検索結果が全く出ないということもある。これは、キーワードの設定仕方によっても大きく左右されるため、入力するキーワードを適切に選ぶことが必要となってくる。

【0005】一般に、この種の検索方法においては、どのようなキーワードを入力したらよいかをユーザ自身で考える必要がある。たとえば、あるキーワードを入力したとき、所望とする文書が得られないような場合には、次に、違うキーワードを入力して検索してみるといった試行錯誤的な検索を行う必要がある。このように試行錯誤的にキーワードを入力して検索を行っても、ユーザが本当に必要としている情報が得られるとは限らない。

【0006】また、情報がある程度絞り込む方法として、複数のキーワードを入力するという方法もあるが、ユーザ自身が何を取り出したらよいのかが明確にわかっていないような場合には、複数のキーワードを設定することは難しいし、また、複数のキーワードによって取り出される情報は、内容が絞り込まれ過ぎることもあり、かえって、所望とする情報を得にくくしてしまう場合もある。。

【0007】ユーザはデータベースに保存されている内容全体を見ることは通常では不可能であるため、いわゆる情報検索という処理を行うわけである。しかし、現在

の情報検索は、データベース内にどのような情報が入っているかが殆どわからない状態で検索を行うために、何をキーワードとしたら最も適切な情報が取り出されるのかがわからないのが実情である。さらに、ユーザ自身、検索すべき情報に対して詳しい知識が無い状態で、どのような情報を得たらよいのか判断できないまま情報検索を行う場合もある。

【0008】このような状況での情報検索を行う場合、従来のように、ユーザの入力したキーワードに基づいて検索を行う方式では、ユーザの所望とする情報を短時間で的確に得ることはできなかった。

【0009】そこで本発明は、データベースの内容を段階的にアウトラインを示しながら表示し、ユーザはその表示を見て選択操作を行うことで、ユーザの必要とする情報を段階的に具体化していくことができるようにし、ユーザ自身がキーワードを考える必要がなく、また、ユーザ自身、検索すべき情報に対して詳しい知識が無い状態で、どのような情報を得たらよいのか判断できないまま情報検索を行う場合でも最終的にユーザの要求する情報を効率よく得ることができる情報検索方法および情報検索装置を提供することを目的としている。

【0010】

【課題を解決するための手段】本発明の情報検索方法において、請求項1の発明では、文書群に属するそれぞれの文書を解析し、それぞれの文書から特徴要素を抽出し、その特徴要素とその特徴要素を含む文書との関係を表す特徴テーブルを作成して、その特徴テーブルに基づいて文書群を複数のクラスタに分類して表示し、ユーザからのクラスタ選択指示を受けたとき、その選択されたクラスタに属する文書に関する内容の表示指示または再分類指示を受け付けて、内容表示指示の場合は、当該クラスタに属する文書に関する表示を行い、再分類指示の場合は、当該クラスタに属する文書のみで前記特徴テーブルを再構成してその再構成された特徴テーブルに基づいてクラスタに分類して表示することを特徴としている。

【0011】また、請求項2の発明は、請求項1の発明において、前記特徴テーブルに基づいて文書群を複数のクラスタに分類する処理は、それぞれの文書内に存在するそれぞれの特徴要素の出現頻度などの統計的な情報に基づいて複数のクラスタに分類するようにしている。

【0012】さらに、請求項3の発明は、請求項2の発明において、前記分類された複数のクラスタは、少なくとも、それぞれのクラスタごとにそのクラスタを代表する特徴要素と、その特徴要素を一定以上含む文書数とをデータとして有している。

【0013】また、本発明の情報検索装置において、請求項4の発明では、文書群を記憶する文書群記憶部と、この文書群記憶部に記憶されているそれぞれの文書を解析する文解析部と、この文解析部による解析結果からそ

れぞれの文書に対する特徴要素を抽出し、その特徴要素とその特徴要素を含む文書との関係を表す特徴テーブルを作成する特徴テーブル作成部と、前記特徴テーブルの内容に基づいて文書群を複数のクラスタに分類する文書分類部と、この文書分類部により分類された内容を記憶する分類結果記憶部と、この分類結果記憶部の内容を読み出して複数のクラスタを表示させる制御を行うとともに、ユーザからのクラスタ選択指示を受けたとき、その選択されたクラスタに属する文書に関する内容の表示指示または再分類指示を受け付けて、内容表示指示の場合は、当該クラスタに属する文書を表示させる制御を行い、再分類指示の場合は、当該クラスタに属する文書のみで前記特徴テーブルを再構成させる制御を行う表示制御部とを有することを特徴としている。

【0014】また、請求項5の発明は、請求項4の発明において、前記文書分類部が行う特徴テーブルに基づいてそれぞれの文書を複数のクラスタに分類する処理は、それぞれの文書内に存在するそれぞれの特徴要素の出現頻度などの統計的な情報に基づいて複数のクラスタに分類するようにしている。

【0015】さらに、請求項6の発明は、請求項5の発明において、前記分類された複数のクラスタは、少なくとも、それぞれのクラスタごとにそのクラスタを代表する特徴要素と、その特徴要素を一定以上含む文書数とをデータとして有している。本発明は、それぞれの文書内に存在するそれぞれの特徴要素の出現頻度などの統計的な情報に基づいて複数のクラスタに分類し、分類された複数のクラスタは、少なくとも、それぞれのクラスタごとにそのクラスタを代表する特徴要素と、その特徴要素を一定以上含む文書数とをデータとして有するようにし、これをユーザに表示するようにしている。

【0016】これにより、ユーザは、表示された複数のクラスタのそれぞれの特徴要素をキーワードとして捉えることができ、それぞれのクラスタごとの特徴要素から、データベース内の概要を知ることができる。したがって、データベース内にどのような情報があるかが全くわからない状態で情報検索を行う場合でも、表示された複数のクラスタについて、ユーザが所望とするクラスタを選択し、かつ、選択されたクラスタ内の文書数が多すぎる場合には、再分類を要求するという処理を段階的かつ対話的に行うことによって、文書数の絞り込みが行え、絞り込まれた状態から最終的にユーザがクラスタを選択し、結果表示要求を行うことで、ユーザの所望とする情報を得ることができる。

【0017】このように、本発明では、情報検索を行うに際して、表示されるクラスタごとの特徴要素を見てクラスタを選択する操作と、再分類結果表示かを状況に応じて選択する操作を行えばよく、その操作過程でユーザの必要とする情報がどれであるかを段々と具体化して行くことができ、最終的にユーザの要求する情報を的確

に取り出すことができる。また、このような情報検索処理を行う過程において、入力すべきキーワードをユーザ自身が考える必要がなく、また、入力操作が選択操作で済むので、検索操作がきわめて容易なものとなる。

【0018】

【発明の実施の形態】以下、本発明の実施の形態を図面を参照して説明する。

【0019】図1は本発明を実現するための装置構成を示す図であり、文書群記憶部11、文解析部12、特徴テーブル作成部13、文書分類部14、分類結果記憶部15、表示制御部16、表示部17などから構成されている。

【0020】文書群記憶部11は、ある文書群に含まれる多数の文書をデータベースとして記憶するものである。

【0021】たとえば、図2に示されるように、ある文書群として、「人工知能に関する論文群」があるとすると、その「人工知能に関する論文群」に属する論文として、たとえば、「エキスパートシステムに関する論文」、「自然言語処理に関する論文」、「ニューラルネットワークに関する論文」、・・・などがあり、さらに、たとえば、「エキスパートシステムに関する論文」の中には、「工場制御に関する論文」、「市場に関する論文」というように、ある1つの文書群には、多数の文書が存在している。

【0022】文解析部12は、文書群記憶部11に記憶されているある文書群のそれぞれの文書データを基に、それぞれの文書の形態素解析を行い特徴要素としての単語を抽出する。

【0023】特徴テーブル作成部13は、特徴要素抽出部131、特徴抽出部132、特徴テーブル133から構成される。特徴要素抽出部131は、文解析部12で形態素解析されて抽出された特徴要素（単語）を基に、それぞれの文書の中に存在する特徴的な単語を抽出する。特徴抽出部132は特徴要素抽出部131から抽出された特徴要素に基づいて、たとえば、それぞれの文書において、どの特徴要素が何回出現したかをカウントする。そして、特徴要素抽出部131から抽出された特徴要素と、特徴抽出部132でカウントされた数とにより、図3のような特徴テーブル133が作成される。

【0024】図3に示される特徴テーブル133の例は、特徴要素としては、「エキスパート」、「システム」、「エキスパート・システム」、「自然」、「言語」、「自然・言語」が示されている。そして、「エキスパート」という特徴要素は、文書Aには4回、文書Bには0回、文書Cには6回、文書Dには3回出現しており、また、「システム」という特徴要素は、文書Aには4回、文書Bには0回、文書Cには8回、文書Dには5回出現しているというように、それぞれの特徴要素がそれぞれの文書にどのくらい出現しているかが示されて

いる。

【0025】この特徴テーブル133の内容によれば、文書Aは、「エキスパート」や「システム」と言った特徴要素が多く出現し、文書Bは「自然」、「言語」、「自然・言語」といった特徴要素が多く出現し、文書Cは「エキスパート」、「システム」、「エキスパート・システム」といった特徴要素が多く出現し、また、文書Dは「エキスパート」、「システム」、「エキスパート・システム」、「自然」、「言語」、「自然・言語」などの特徴要素がどれも多く出現していることがわかる。

【0026】文書分類部14は、このような内容の特徴テーブル133を参照し、それぞれの文書内に存在するそれぞれの特徴要素の出現頻度などの統計的な情報に基づいて複数のクラスタに分類する。

【0027】たとえば、1番目のクラスタとしては、そのクラスタを代表する特徴要素が「エキスパート・システム」であり、その「エキスパート・システム」という特徴要素を一定以上含む文書数は「2」であり、2番目のクラスタとしては、特徴要素が「自然・言語」であり、その「自然・言語」という特徴要素を一定以上含む文書数は「3」であるというように分類される。また、文書分類部14は、このような特徴要素とその特徴要素を一定以上含む文書数の分類を行うとともに、それがどの文書であるかについての対応付けも行う。たとえば、文書数「2」に対応する文書は文書Aと文書Cであるという文書数と文書名の対応付けも行う。このような分類結果は分類結果記憶部15に格納される。

【0028】表示制御部16は、分類結果記憶部15の内容を分類結果として表示部17に表示する制御を行うとともに、ユーザから検索結果表示の指示があったときは、分類結果記憶部15の内容と前記文書群記憶部11の内容に基づいて検索結果内容としての表示データを構成し、その表示データを表示部17に表示する制御を行う。

【0029】図4は表示部17に表示された分類結果の一例を示すもので、この例では1つの画面上には、クラスタとしてたとえば分類1から分類10まで、10個のクラスタを表示する。なお、ここでは、1画面に10個単位としたが、これは適当な数を設定できるものであり、また、クラスタ数が多い場合は、10個ずつに分けてページ切替えで表示するようにすることも可能である。

【0030】この図4の例では、分類1の特徴要素は「エキスパート・システム」であり、その文書数は「2」、分類2としては、特徴要素が「自然・言語」であり、その文書数は「3」であるというように表示されている。このように、各クラスタ毎にそのクラスタを代表する特徴要素とその特徴要素が一定以上存在する文書数が表示される。また、その表示部17には「結果表示」と「再分類」といったユーザの指示を入力するため

のユーザ指示部21、22が表示される。

【0031】ユーザはこのような表示内容を見て、ユーザ自身の要求している情報が、たとえば、分類1の内容（「エキスパート・システム」）に関係するものではないかと判断した場合は、その分類1の行部分R1をマウスなどでクリックしたのち、「結果表示」のユーザ指示部21をクリックする。

【0032】これにより、表示制御部16は、選択されたクラスタ（分類1）に属する文書を文書群記憶部11から読み出して、その文書内容を表示部17に表示する。この例では、選択されたクラスタ（分類1）に属する文書数は「2」であり、その文書名が文書Aと文書Cであることがわかるから、表示制御部16は、ユーザからの結果表示要求を受けると、ユーザの選択したクラスタに属する文書（文書Aと文書C）を文書群記憶部11から読み出して、その内容を表示する。

【0033】なお、この文書内容の表示の仕方としては、分類1に属するすべての文書（この場合文書Aと文書C）の内容をそのまますべてを表示させるようにしてもよいが、たとえば、文書が論文である場合には概要を表す部分のみを表示するようにしてもよく、あるいは、文書名と文書サイズなどのみを表示するようにしてもよく、その表示の仕方については種々考えられる。

【0034】一方、ユーザが図4に示すような表示内容を見て、ユーザ自身の要求している情報が、たとえば、分類1の内容（「エキスパート・システム」）に関するものではないかと判断したものの、「エキスパート・システム」という表示内容だけでは、ユーザ自身の要求する情報として具体化されていないと判断した場合、つまり、もう少し細分類化された内容が必要であると考えたときは、分類1の行部分R1をマウスなどでクリックしたのち、「再分類」のユーザ指示部22をクリックする。

【0035】このように、分類1の行部分R1がクリックされたのち、「再分類」のユーザ指示部22がクリックされると、分類1に属する文書のみで再分類処理される。この図4に示す例では、分類1に属する文書数は「2」である。したがって、この2つの文書のみを用いて、それらの文書に存在する特徴要素に基づいて特徴テーブルを再構成する。つまり、この例で考えると、分類1に属する文書は文書Aと文書Cであるから、これらの文書Aと文書Cとで新たな特徴テーブル133が作成されることになる。そして、新たに作成された特徴テーブルを参照して、文書分類部14が文書群の分類を行い、この文書Aと文書Cのみについて分類されされた内容が表示部17に表示される。

【0036】このようにして、分類対象の文書が絞られた状態で、再分類された結果は、再分類前の特徴要素がある程度はそのまま出てくるが、新たな分類対象の文書間で見た場合、ある文書に特有の特徴要素が、分類結果

として出てくる場合もある。たとえば、分類対象の文書を文書Aと文書Bとしたとき、文書Cでは「工場制御」という特徴要素の出現頻度が高いが、この「工場制御」という特徴要素は文書Aでは殆ど出現しないという場合は、1つのクラスタとして、特徴要素が「工場制御」でその「工場要素」を含む文書数が「1」というような分類結果が表示されることになる。この表示例を図5に示す。図5では、分類3のクラスタにおいて、特徴要素が「工場制御」でその「工場制御」を一定以上含む文書数が「1」というように表示されている。

【0037】そして、ユーザがその再分類された表示結果を見て、ユーザの要求する情報が「工場制御」に関する内容に近いと判断した場合には、図5における分類3の行部分R2をクリックし、かつ、「結果表示」のユーザ指示部21をクリックすると、文書Cの内容が表示される。なお、この表示についても前記したように、対象となる文書の内容をそのまま表示させるようにしてもよいが、たとえば、文書が論文である場合には概要を表す部分のみを表示するようにしてもよく、あるいは、文書名と文書サイズなどのみを表示するようにしてもよく、その表示内容については種々考えられる。

【0038】なお、以上の例は、説明を容易なものとするために、図4の段階で分類される文書数を「2」というようなきわめて少ない数で説明したが、実際には、図4の段階では、それぞれの分類における文書数は数百というような数となることもある。したがって、ユーザがたとえば、分類1のクラスタを選択し、かつ、再分類を指示すると、その分類1に属する数百の文書での再分類がなされ、その再分類された内容として、分類1における特徴要素とその文書数、分類2における特徴要素とその文書数、分類3における特徴要素とその文書数というように、それぞれの分類番号ごとにその特徴要素とその特徴要素を含む文書数が表示部17に表示される。

【0039】そして、ユーザが、その再分類されたそれぞれの分類番号に対する特徴要素を見て、ある分類番号のクラスタを選択し、かつ、再分類を要求すると、今度は、その選択した分類番号に属する文書だけの再分類がなされる。たとえば、ユーザが、分類3のクラスタを選択し、かつ再分類を指示すると、その分類3に属する文書数での再分類がなされ、前記同様に、分類1における特徴要素とその文書数、分類2における特徴要素とその文書数、分類3における特徴要素とその文書数というように、それぞれの分類番号ごとにその特徴要素とその特徴要素を含む文書数が表示される。

【0040】このような処理が繰り返行われることにより、対象文書数が段階的に絞り込まれて行く。そして、文書数が絞り込まれた状態で、ユーザは表示された特徴要素を見て、最も適当と思われる特徴要素が表示されているクラスタ部分をクリックしたのち、「結果表示」のユーザ指示部21をクリックする。

【0041】これにより、たとえば、最終的な段階でユーザの選択した特徴要素を含む文書数が「2」であれば、その特徴要素を含む2つの文書の内容が表示されることになる。なお、この結果表示処理は、分類結果記憶部15に記憶されている最新の分類結果内容と、文書群記憶部11の内容を基に、対応する文書名が読み出され、検索結果となる表示データを構成して、その表示データを表示部17に表示することにより行う。

【0042】このように、本発明では、表示部17に表示される分類結果(図4参照)における分類番号に対応する特徴要素がいわばキーワードとなるものである。

【0043】したがって、ユーザは自分の要求する情報について、何をキーワードとしてよいかかわからないような場合であっても、システム側で、ユーザの要求を具体化するための指標となる特徴要素を画面上に分類番号対応に表示し、さらに、その特徴要素を一定以上含む文書の数を表示するので、ユーザはデータベースの概要を知ることができ、また、ユーザ自身がキーワードを考える必要がなく、その表示内容を見て、選択するという対話的な操作が可能となる。

【0044】そして、必要に応じて何段階かの再分類操作を経て文書数が絞り込まれたところで、最も適当と思われる特徴要素を選択してその結果表示を行うというような検索処理を行うことで、必要とする文書を的確に取り出すことができる。

【0045】図6は以上説明したこの実施の形態の処理手順をフローチャートである。図6において、文書分類部14が特徴テーブル133を参照して文書群の分類を行い(ステップs1)、その分類結果を表示する(ステップs2)。この分類結果の一例としては、たとえば、図4で示すような内容である。そして、ユーザがその表示を見て、「結果表示」か「再分類」かの入力を行うと、そのユーザ指示入力を受け付け(ステップs3)、ユーザの指示が結果表示か再分類かを判断する(ステップs4)。ユーザ指示が結果表示である場合には、文書群記憶部11と分類結果記憶部15のそれぞれの内容から表示データを構成して(ステップs5)、その表示データを表示部17に表示する(ステップs6)。

【0046】一方、ステップs4において、ユーザの指示が再分類である場合には、選択されたクラスタ内の文書のみで特徴テーブル133を再構成する(ステップs7)。そして、ステップs1に処理が戻り、再構成された特徴テーブルを用いて文書群を分類し、以下前記ステップs2～s7の処理を行う。なお、指示されたクラスタ内の文書のみで特徴テーブル133を再構成する処理は、指示されたクラスタの文書が前述したように、たとえば、文書Aと文書Cであるとすれば、この文書Aと文書Cのみを用いて、それぞれの文書から抽出された特徴要素に基づいて特徴テーブル133を再構成する処理である。

【0047】以上説明したように、この実施の形態によれば、特徴要素としての単語の出現頻度などの統計的情報によって、文書群を自動的に、所定の数のクラスタに分類するとともに、それぞれのクラスタを代表する特徴要素とその特徴要素を一定以上含む文書の数を抽出し、それらをユーザに見せることによって、ユーザは、データベースの概略を知ることができる。そして、このようにして分類されたクラスタについて、ユーザはそれぞれのクラスタごとの特徴要素をキーワードとして捉えて、最も適当と思われる特徴要素の存在するクラスタを選択する。

【0048】このとき、選択したクラスタ内の文書数が多すぎる場合には、選択されたクラスタ内の文書のみで再分類して表示することが可能で、この再分類された表示を見て、その中で、最も適当と思われる特徴要素の存在するクラスタを選択するという操作を段階的かつ対話的に行う。これにより、文書数の絞り込みが行え、絞り込まれた状態から最終的にユーザがクラスタを選択し、結果表示要求を行うことで、ユーザの所望とする情報を得ることができる。

【0049】なお、以上説明した実施の形態は、本発明の好適な実施の形態の一例であるが、本発明はこれに限定されるものではなく、本発明の用紙を逸脱しない範囲で種々変形実施可能となるものである。また、本発明の処理を行う処理プログラムは、フロッピーディスク、光ディスク、ハードディスクなどの記憶媒体に記憶しておくことができ、本発明は、それらの記憶媒体をも含むものであり、また、ネットワークからデータを得る形式でもよい。

【0050】

【発明の効果】本発明によれば、それぞれの文書から特徴要素を抽出し、その特徴要素とその特徴要素を含む文書との関係を表す特徴テーブルを作成して、その特徴テーブルに基づいて文書群を複数のクラスタに分類して表示し、ユーザからのクラスタ選択指示を受けたとき、その選択されたクラスタに属する文書に関する内容の表示指示または再分類指示を受け付けて、内容表示指示の場合は、当該クラスタに属する文書に関する表示を行い、再分類指示の場合は、当該クラスタに属する文書のみで前記特徴テーブルを再構成してその再構成された特徴テーブルに基づいてクラスタに分類して表示するようにしている。これにより、ユーザは、表示された複数のクラスタのそれぞれの特徴要素をキーワードとして捉えることができ、それぞれのクラスタごとの特徴要素から、データベース内の概要を知ることができる。したがって、データベース内にどのような情報があるかが全くわからない状態で情報検索を行う場合でも、クラスタごとの特徴要素からデータベース内の概要を知ることができる。そして、分類された複数のクラスタは、少なくとも、それぞれのクラスタごとにそのクラスタを代表する特徴要

素と、その特徴要素を一定以上含む文書数とをデータとして有するので、表示された複数のクラスタについて、ユーザが所望とするクラスタを選択したとき、選択したクラスタ内の文書数が多すぎる場合には、再分類が可能で、この処理を段階的かつ対話的に行うことによって、文書数の絞り込みが行え、絞り込まれた状態からユーザが所望とする情報を得ることができる。

【0051】このように、本発明では、情報検索を行うに際して、表示されるクラスタごとの特徴要素を見てクラスタを選択する操作と、再分類か結果表示かを選択するという対話的でかつ段階的に検索対象を絞って行く操作を行えばよく、その操作過程でユーザの欲する情報がどれであるかを段々と具体化して行くことができ、最終的にユーザの必要とする情報を的確に取り出すことができる。また、このような情報検索処理を行う過程において、入力すべきキーワードをユーザ自身が考える必要がなく、また、入力操作がマウスなどによる選択操作ですむので、検索操作がきわめて容易なものとなる。

【図面の簡単な説明】

【図1】本発明の実施の形態における情報検索装置の構成を示すブロック図。

【図2】本発明の実施の形態に用いられる文書群の例を示す図。

【図3】本発明の実施の形態における特徴テーブルの一例を示す図。

【図4】本発明の実施の形態における複数のクラスタ表示例を示す図。

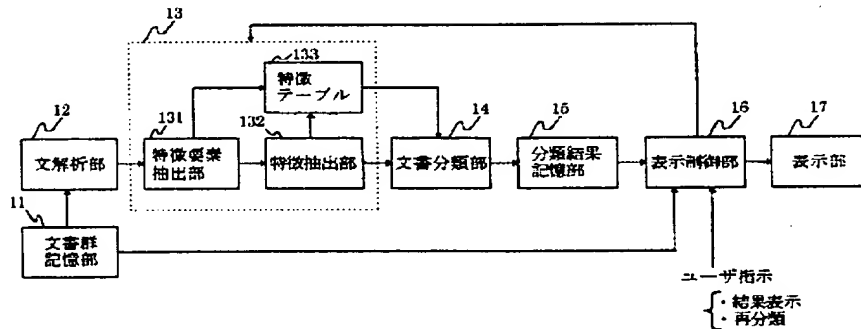
【図5】本発明の実施の形態における再分類指示により再構成された特徴テーブルに基づいてクラスタに分類された表示例を示す図。

【図6】本発明の実施の形態における検索処理手順を説明するフローチャート。

10 【符号の説明】

- 11 文書群記憶部
- 12 文解析部
- 13 特徴テーブル作成部
- 14 文書分類部
- 15 分類結果記憶部
- 16 表示制御部
- 17 表示部
- 21 「結果表示」のユーザ指示部
- 22 「再分類」のユーザ指示部
- 131 特徴要素抽出部
- 132 特徴抽出部
- 133 特徴テーブル

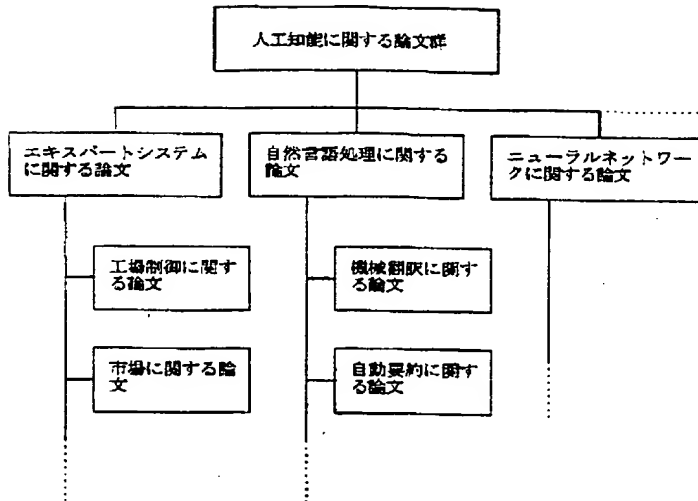
【図1】



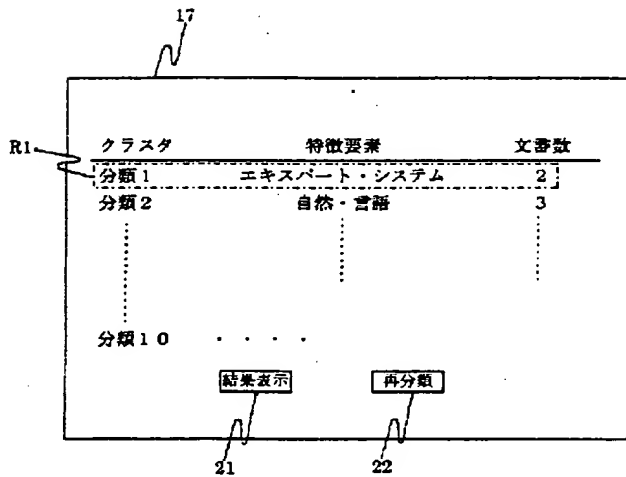
【図3】

特徴要素	文書A	文書B	文書C	文書D	...
エキスパート	4	0	6	3	...
システム	4	0	8	5	...
エキスパート・システム	2	0	4	2	...
自然	0	5	3	3	...
言語	1	6	0	3	...
自然・言語	0	5	0	3	...

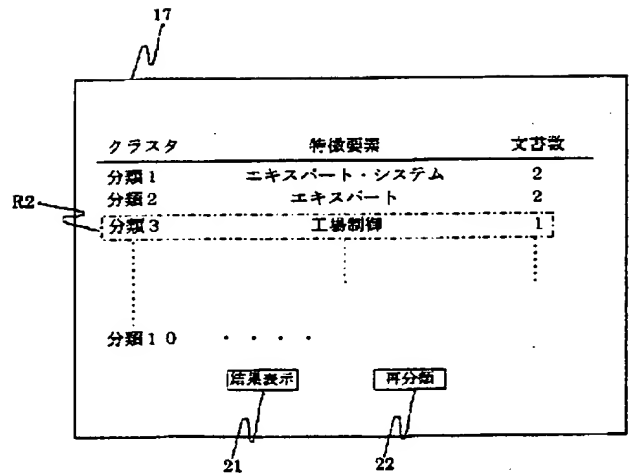
【図2】



【図4】



【図5】



【図6】

